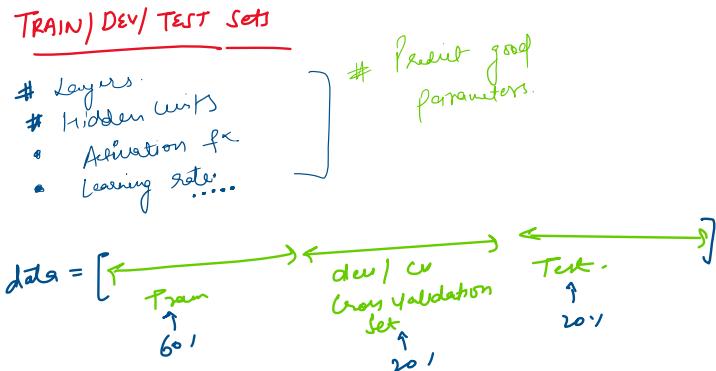


## Setting up your Machine Learning Application.

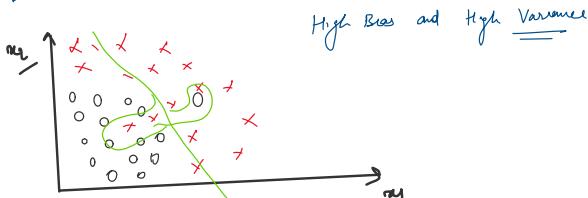


# Not having a test set is okay and only have Dev set.

### Bias / Variance

Train set Error  $\rightarrow 17$ .  
Dev set Error  $\rightarrow 117$ .

Underfitting - High Bias



### Basic Recipe for Machine Learning

- ① High Bias ??  $\rightarrow$  \* Bigger N/w Layer -  
(Training data performance)  
• Layer optimization Algos
- ② High Variance ?  $\rightarrow$  \* More data.  
(Dev set performance)  
• Regularization  
• neural N/w Architectures;

Done!

### bias / Variance Tradeoff

bias ↑ Variance ↓

## Regularizing your Neural N/w

$$\|w^u\|^2 = \sum_{i=1}^m \sum_{j=1}^{n(u)} (w_{ij}^{(u)})^2$$

### # Regularization

$$\text{Logistic Regression} - \min_{(w,b)} J(w,b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2} \|w\|^2$$

L<sub>2</sub> regularization  $\rightarrow w \in \mathbb{R}^{n \times k}, b \in \mathbb{R}$

$$\|w\|^2 = \sum_{j=1}^n w_j^2 = w^T w$$

$$\|w\|_1 = \sum_{j=1}^n |w_j| \quad (\lambda = \text{Regularization parameter})$$

L<sub>2</sub> Regularization

$$\frac{1}{m} \sum_{i=1}^m \|w\|_2^2 = \frac{\lambda}{m} \|w\|_F^2$$

( $\lambda$  = Regularization parameter)  
[Hyperparameter]

Neural N/w

$$J(w^{(0)}, b^{(0)}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|_F^2$$

$$\|w^{(l)}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^{n^{(l+1)}} (w^{(l)}_{ij})^2$$

} "Frobenius norm"  $\| - \|_F$

$$w \in (\mathbb{R}^{n^{(l)}}, \mathbb{R}^{n^{(l+1)}})$$

$$dw = (\text{Calculated from Backprop}) + \frac{\lambda}{m} w$$

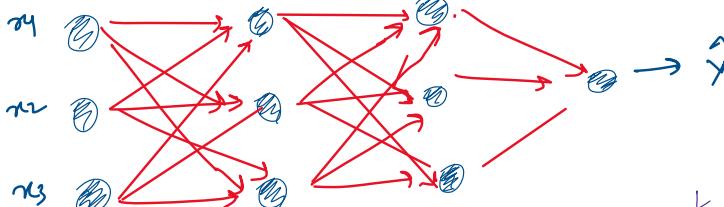
$$\frac{\partial J}{\partial w^{(l)}} = dw^{(l)}$$

$$w^{(l)} = w^{(l)} - \alpha dw^{(l)}$$

# L<sub>2</sub> regularization is also called "weight decay"

$$\begin{aligned} w^{(l)} &:= w^{(l)} - \alpha \left[ (\text{Backprop}) + \frac{\lambda}{m} w^{(l)} \right], \\ &= w^{(l)} - \frac{\alpha \lambda}{m} w^{(l)} - \alpha (\text{Backprop}) \end{aligned}$$

Why Regularization prevents Overfitting?

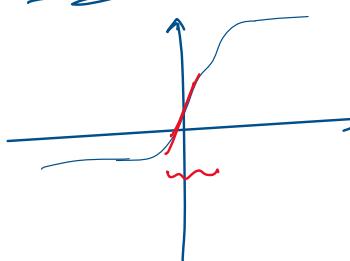


$$J(w^{(0)}, b^{(0)}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|_F^2$$

① (If  $\lambda \gg$  big,  $w \approx 0$ )

② high bias  $\rightarrow$  high variance

tanh(f(x))

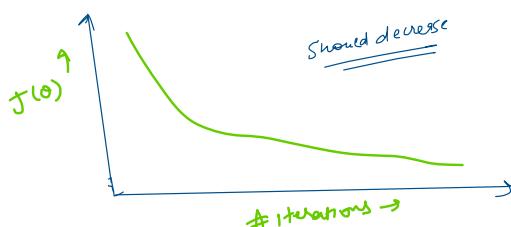


$$g(x) = \tanh(x)$$

$w^{(l)} \downarrow, \lambda \uparrow$   
 $f(x)$  will be linear  
in red range

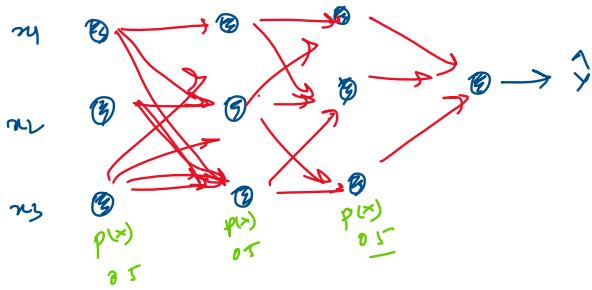
# As if every layer is linear  $N/w$

$$J(\cdot) = \sum L(\hat{y}, y) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|_F^2$$



## Dropout    REGULARIZATION

④ Set some probability to dropout a node in now,



## Implementing Dropout (Inverted dropout)

# Illustrate in layer  $L=3$ .

$\text{ds} = \text{np.random.rand}(\text{a3.shape[0]}, \text{a3.shape[1]}) < \underline{\underline{\text{keep prob.}}}$

$$a_3 = \text{np multiply } (a_1, d_3) \quad (a_3 \neq d_3)$$

$a_3$  / = keep prob = (Inverted dropout technique)

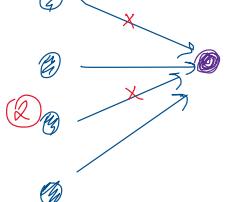
④ (don't dropout at test time.)

④ (don't dropout at test time.)

## Understanding Dropout

(#) Can't rely on one feature, so have to spread weights  
≈ Share weights

~~(A) ④~~ 3 4 5



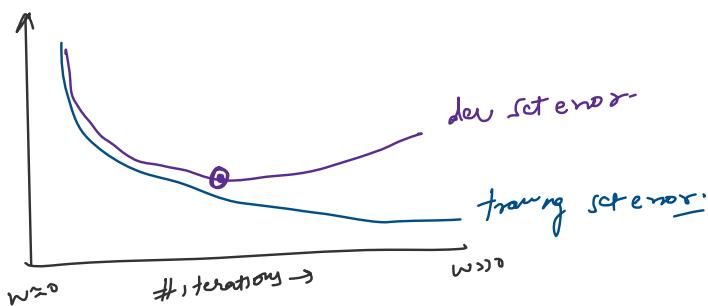
④ Can vary, keep prop according to layer,

$$\# \text{ units} \propto \frac{1}{\text{keep prob value}}$$

## Computer Vision

## Other Regularization Methods

## EARLY STOPPING



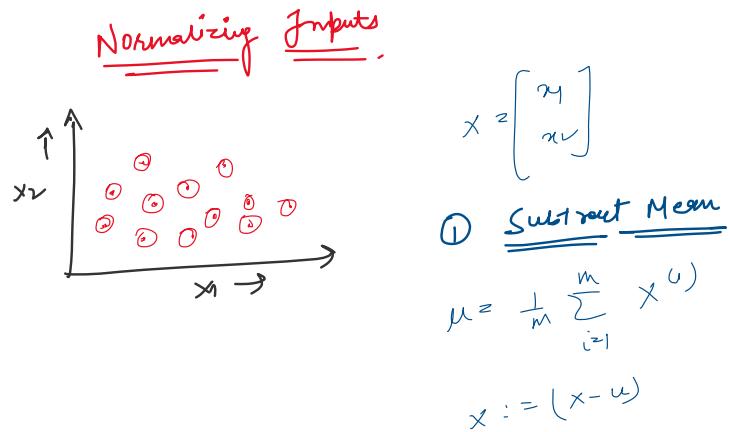
# optimize cost fx or (Gradient descent, LBGK, ...)

# do not overfit  
(regularization, getting more data)

Orthogonalization → Thinking of one task at a time.

- Orthogonalization → Thinking of one task at a time  
 → ① "no orthogonalization in early stopping"
- $L_2$  Regularization

## Setting up your Neural N/W



② Normalize Variances

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m [x^{(i)} - \mu]^2 \quad \text{(Squaring.)}$$

$$x = \frac{x}{\sigma}$$

$$x = \frac{(x - \mu)}{\sigma}$$

③ use same  $\mu$  and  $\sigma^2$  to normalize Test set.

Why NORMALIZE?

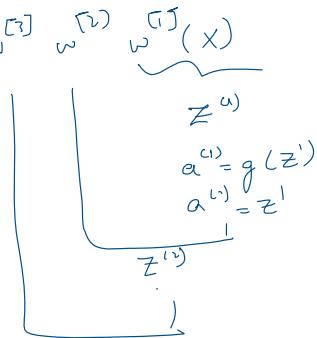
$$J(\omega, b) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)})$$

$$\begin{aligned} w_1 &= 1 && \dots & 1000 \\ w_2 &= 1 && \dots & 10 \end{aligned}$$

Variance | Exploding gradients



$$\hat{y} = w^{[1]} + w^{[2]} * g(z) = \hat{z}$$



$$W^{[l]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix} \quad \text{with } \lambda^{(l)} = 10^5$$

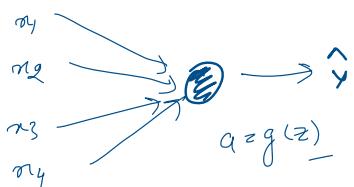
$$\hat{y} = W^{[l]} \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{-1} \times \hat{y}$$

$\hat{y} = (1.5)^L$

- $\rightarrow 10^5$  times number of layers
- $\rightarrow$  if you have deep neural net, value of  $\hat{y}$  will explode.

## WEIGHT INITIALIZATION FOR NEURAL N/W.

### # Single Neuron



$$z = w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

$$b = 0$$

$$n \propto \frac{1}{w_i}$$

$$\text{Var}(w_i) = \frac{1}{n} \quad [n \in \text{No. of input features in neuron}]$$

$$W^{[l]} = np.random.randint(\text{shape}) * np.sqrt\left(\frac{1}{n^{[l-1]}}\right)$$

$2/n$  for ReLU ( $f(x)$ )

$$g^{(l)}(z) \in \text{ReLU}$$

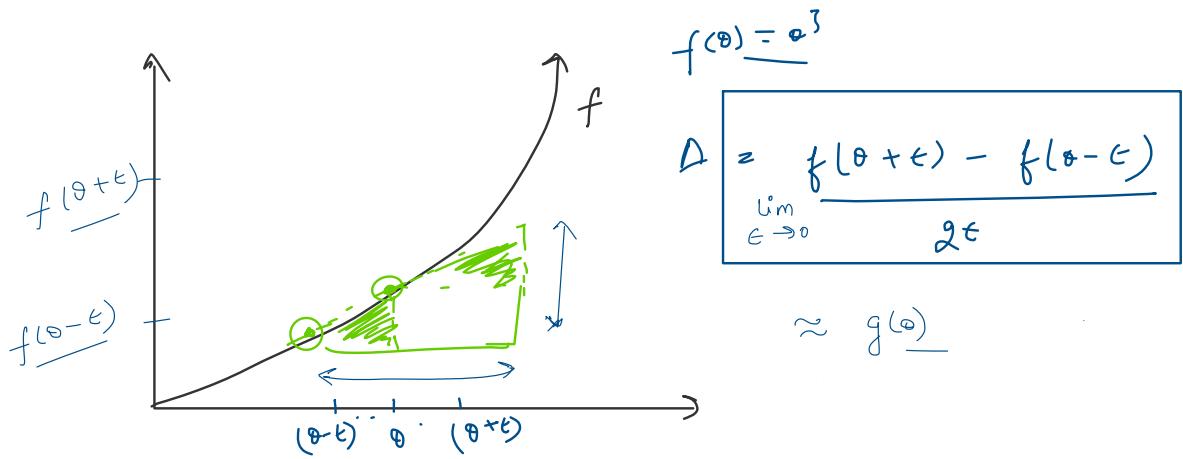
### Other Variants

$$\tanh(z) = \sqrt{1/n^{[l-1]}} \quad \left. \begin{array}{l} \text{xavier} \\ \text{initialization} \end{array} \right\}$$

## Numerical Approximation OF GRADIENTS

~ Backprop (Gradient checking) -

$$\uparrow \quad \uparrow \quad f(\theta) = 0^3$$



$$f(\theta) = \underline{\theta^3}$$

$$\Delta = \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \underset{\epsilon \rightarrow 0}{\approx} g(\theta)$$

### Gradient Checking.

gradient checking for Neural N/W  $\rightarrow$

Take  $w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)}$  & reshape to a big vector,  $\underline{\theta}$ ,  
Concatenate

Take  $d\theta^{(1)}, d\theta^{(2)}, d\theta^{(3)}, \dots, d\theta^{(L)}$  & reshape into a big vector,  $d\underline{\theta}$   
concatenate  $\rightarrow$

grad check

for each  $i \rightarrow$

$$d\theta_{approx}^{(i)} := \frac{J(\theta_1, \theta_2, \dots, \theta_i + \epsilon) - J(\theta_1, \theta_2, \dots, \theta_i - \epsilon)}{2\epsilon}$$

$$\approx d\theta^{(i)}$$

$$= \frac{\partial J}{\partial \theta_i}$$

$$\partial \theta_{approx} \approx \partial \theta$$

$$\underline{\underline{\text{check}}} = \frac{\|\partial \theta_{approx} - \partial \theta\|_2}{\|\partial \theta_{approx}\| + \|\partial \theta\|_2}$$

$$\epsilon \approx 10^{-7}$$

### gradient Checking Implementation

- ④ don't use grad check in training only to debug,

\* if algorithm fails in grad check; look at the components & try to find the bug.

$$\partial b^{[l]}, \partial w^{[l]}$$

\* Remember regularization

$$J(\theta) = \frac{1}{m} \sum L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_2^2$$

$$\partial \theta = \text{grad of } J \text{ wrt. } \theta$$

\* grad check doesn't work with dropout

\* run at random initialization; perhaps after some training

Mini batch gradient descent

$$\mathbf{X} = [\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \mathbf{x}^{(3)}, \dots, \mathbf{x}^{(m)}] \\ (n \times m)$$

$$\mathbf{Y} = [y^{(1)}, y^{(2)}, y^{(3)}, \dots, y^{(m)}] \\ (1, m)$$

if training examples are very large, say 10 million.

① will have to traverse (10 million) training set before you take one little step of gradient descent.

Mini-Batch

Now, One mini batch consists of

$$\text{if } X=50 \text{ (say)} \text{ & } n=5$$

$$X^{\{1\}} = \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & x^{(4)} & \dots & x^{(5)} \end{bmatrix}$$

$$X^{\{1\}} = \dots$$

$$\frac{\text{training-examples}}{\text{n. of mini-batches}} = n(\text{mb}).$$

→ and same for Y.

# Mini-Batch t: $X^{t+3}, Y^{t+3}$	① $X_{\text{dim}} = n_x, n(\text{mb})$
	② $Y_{\text{dim}} = 1, n(\text{mb})$

Mini-Batch gradient descent

for  $t = 1 \dots n(\text{mb})$ :

// perform a step of gradient descent using  $X^{\{t\}}, Y^{\{t\}}$

① Perform forward prop.  $\Rightarrow$

$$Z^{[0]} = W^{[0]} X^{[0]} + b^{[0]}$$

$$A^{[0]} = g^{[0]}(Z^{[0]})$$

$$A^{[L]} = g^{[L]}(Z^{[L]})$$

② Compute loss

$$J = \frac{1}{m} \sum J(\rightarrow) + \text{regularization term}$$

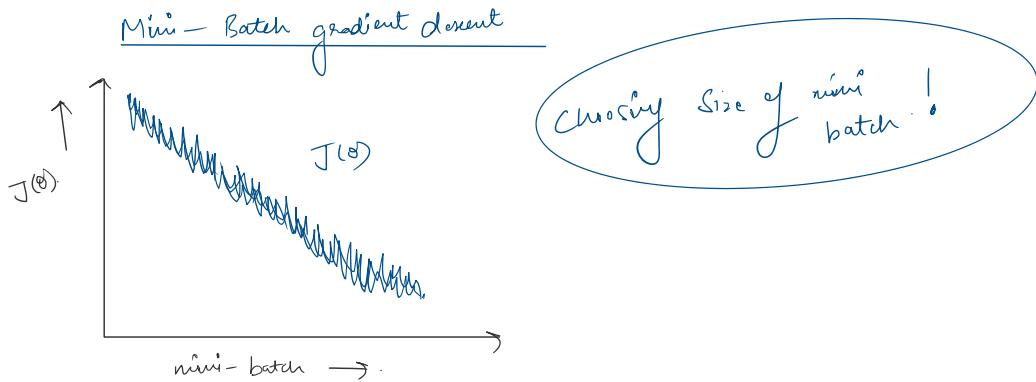
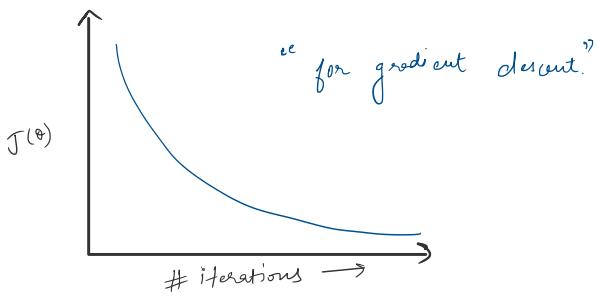
$$W^{[L]} := W^{[L]} - \alpha \partial W^{[L]} \quad \& \text{ similarly,}$$

$$b^{[L]} := b^{[L]} - \alpha \partial b^{[L]}$$

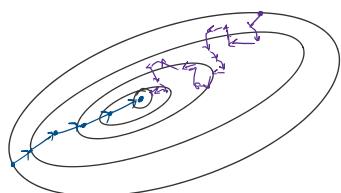
"Also called 1 epoch / 1 pass"

L2. UNDERSTANDING MINI BATCH GRADIENT DESCENT.

## Ld. Understanding Mini-Batch Gradient Descent.



- ① if mini batch size =  $m$ , then : Batch gradient descent.
- ② if (mini batch size = 1) : Stochastic gradient descent.  
 $(x^{(1)}, y^{(1)}, x^{(2)}, y^{(2)}, \dots)$



[In practice mini batch size is between 1 and  $m$ .]

### Batch gradient descent

mini-batch size :  $m$

↓  
Too long per iteration

### Stochastic gradient descent

batch size : 1

↓  
loop speeding up from vectorization.

④ mini - batch gradient descent

### Choosing size of mini-batch

if Small training set : Use batch gradient descent.

if Small training set : Use batch gradient descent.  
( $m \leq 2000$ )

Typical mini-batch sizes :  $2^k$  (multiple of 2.)

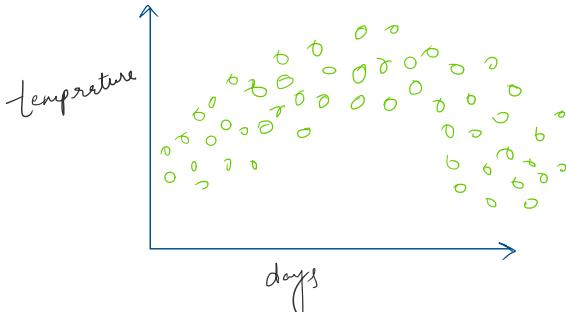
Make sure mini-batch fits in CPU / GPU memory.

$X^{(t)}$ ,  $y^{(t)}$  → Hyperparameter.

## L3 EXPONENTIALLY WEIGHTED AVERAGE

Temperature in New-delhi →

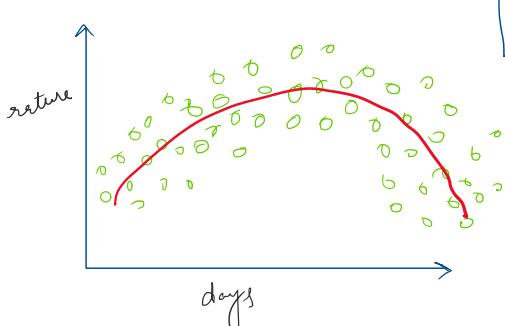
$$\begin{aligned} \theta_1 &= 32^\circ C \\ \theta_2 &= 31^\circ C \\ \theta_3 &= 29^\circ C \\ &\vdots \\ \theta_{190} &= 28^\circ C \end{aligned}$$



$$\begin{aligned} V_0 &= 0 \\ V_1 &= 0.9V_0 + 0.1\theta_1 \\ V_2 &= 0.9V_1 + 0.1\theta_2 \\ V_3 &= 0.9V_2 + 0.1\theta_3. \end{aligned}$$

$$V_t = 0.9V_{(t-1)} + 0.1\theta_t$$

→ Exponentially weighted average.  
→ Moving Average.



$$V_t = \beta V_{(t-1)} + (1-\beta)C$$

Exponentially weighted moving average.

generalize formula.

$$\beta = 0.9$$

$$V_t \stackrel{(avg)}{=} \frac{1}{(1-\beta)} \text{ days } \text{ temperature.}$$

L4

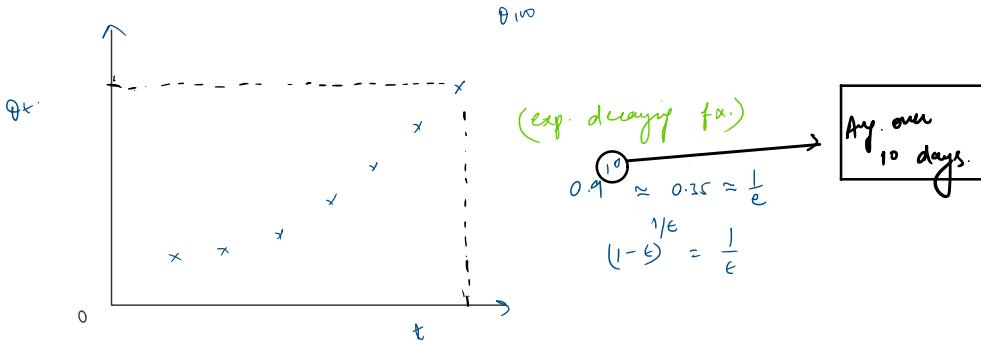
## UNDERSTANDING EXPONENTIALLY WEIGHTED AVG.

$$V_t = \beta V_{(t-1)} + (1-\beta) \theta_t$$

$$V_{100} = 0.1 \theta_{100} + 0.9 V_{99} \quad \beta = 0.9$$

$$= 0.1 \theta_{100} + 0.1 * 0.9 (\theta_{99}) + 0.1 (0.9)^2 \theta_{98} + \dots$$

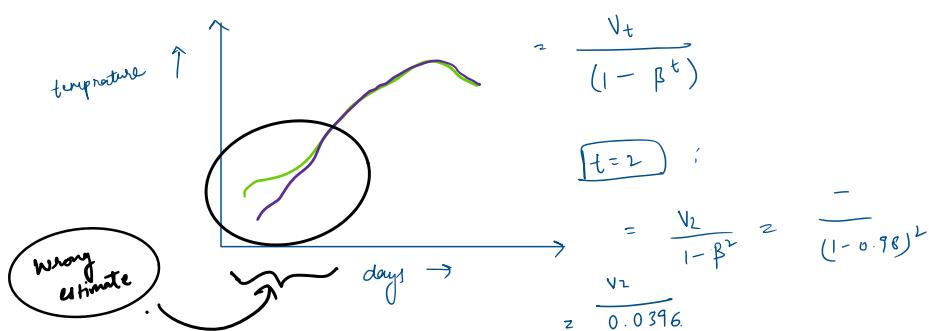
- exponentially decaying function.



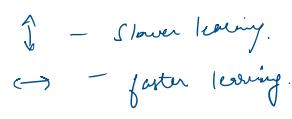
$\cdot V_0 = 0$ $\cdot V_t = \beta V_{(t-1)} + (1-\beta) \theta_t$	$V = 0$ $V_0 = \beta V + (1-\beta) \theta_1$ $V_1 = \beta V_0 + (1-\beta) \theta_2$
---	---

## L5 BIAS CORRECTION IN EXP. WEIGHTED AVG.

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t \quad ; \quad \beta = 0.98$$



## L6 GRADIENT DESCENT WITH MOMENTUM



# Momentum :-

On every iteration  $t$ , calculate  $\nabla w, \nabla b$  on current batch

$\nabla w, \nabla b$  on current batch  $\rightarrow$  friction

$$\nabla w = \beta \nabla w + (1-\beta) \nabla w \quad ; \quad \text{Acceleration}$$

friction

$\delta w, \delta b$  on current batch

Velocity

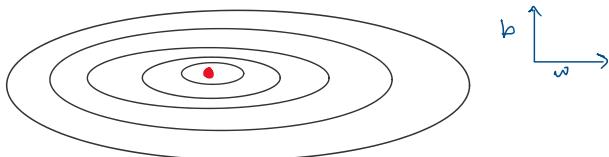
$$\begin{aligned} V_{\delta w} &= \beta V_{\delta w} + (1-\beta) \delta w \\ V_{\delta b} &= \beta V_{\delta b} + (1-\beta) \delta b \end{aligned} \quad \left. \begin{array}{l} \text{Acceleration} \\ \vdots \end{array} \right.$$

$$w := w - \alpha V_{\delta w}$$

$$b := b - \alpha V_{\delta b}$$

Hyperparameters :-  $\alpha, \beta$        $\boxed{\beta = 0.9}$ .

## L7. RMS Prop (Root mean square prop)



On every iteration, t :-

Compute  $\delta w, \delta b$  for mini-batch.

$$S_{\delta w} = \beta S_{\delta w} + (1-\beta) \delta w^2$$

element wise square  $\Rightarrow$   $p$  is hyperparam for momentum

$$S_{\delta b} = \beta S_{\delta b} + (1-\beta) \delta b^2$$

$$\boxed{S_{\delta w} \approx}$$

to make sure it's tve, adding small value  $\epsilon$ .  
 $\epsilon \approx 10^{-8}$

$$w := w - \frac{\alpha \delta w}{\sqrt{S_{\delta w}} + \epsilon}$$

①  $w \rightarrow$  : fast.  
 $b \uparrow$  : slow down.

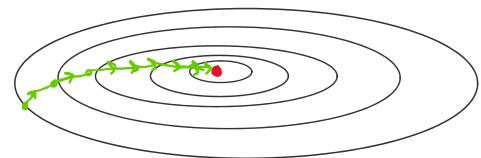
$$b := b - \frac{\alpha \delta b}{\sqrt{S_{\delta b}} + \epsilon}$$

$S_{\delta w}$  = relatively small.

$S_{\delta b}$  = relatively large.

\* RMS prop.

→ now, can also use large learning rate,  $\alpha$ .



## L8. ADAM'S OPTIMIZATION ALGORITHM

$$\{ \text{grad (Momentum)} + \text{Rmsprop.} \}$$

$$V_{\delta w} = 0, \quad S_{\delta w} = 0 ; \quad V_{\delta b} = 0, \quad S_{\delta b} = 0$$

On every iteration t :

\* Compute  $\delta w, \delta b$  using current mini-batch.

$$\begin{aligned} V_{\delta w} &= \beta_1 V_{\delta w} + (1-\beta_1) \delta w \\ V_{\delta b} &= \beta_1 V_{\delta b} + (1-\beta_1) \delta b \end{aligned} \quad \left. \begin{array}{l} \text{• Momentum (Hyp: } \beta_1 \text{)} \\ \vdots \end{array} \right.$$

$$\begin{aligned} S_{\delta w} &= \beta_2 S_{\delta w} + (1-\beta_2) \delta w^2 \\ S_{\delta b} &= \beta_2 S_{\delta b} + (1-\beta_2) \delta b^2 \end{aligned} \quad \left. \begin{array}{l} \text{• rms prop.} \\ (\text{Hyp: } \beta_2) \end{array} \right.$$

### Bias correction

$$* \quad \begin{aligned} V_{dw}^{\text{corrected}} &= V_{dw} / (1 - \beta_1^t) \\ V_{db}^{\text{corrected}} &= V_{db} / (1 - \beta_1^t) \end{aligned} \quad \left. \right\} \cdot \text{Momentum}$$

$$\left. \begin{aligned} S_{dw}^{\text{corrected}} &= S_{dw} / (1 - \beta_2^t) \\ S_{db}^{\text{corrected}} &= S_{db} / (1 - \beta_2^t) \end{aligned} \right\} \cdot \text{Rms Prop.}$$

$$* \quad w_i = w_i - \alpha \left[ \frac{V_{dw}^{\text{corrected}}}{(\sqrt{S_{dw}^{\text{corrected}}} + \epsilon)} \right]$$

$$* \quad b_i = b_i - \alpha \left[ \frac{V_{db}^{\text{corrected}}}{(\sqrt{S_{db}^{\text{corrected}}} + \epsilon)} \right]$$

### hyperparameters choice :-

$\alpha$  = Needs to be tuned.

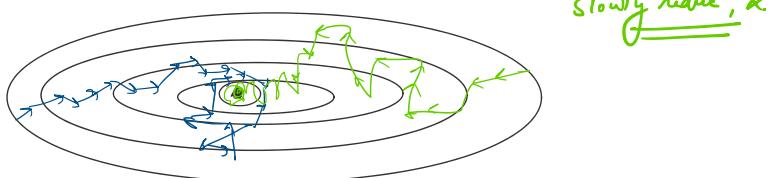
$\beta_1 = 0.9$  (dw)

$\beta_2 = 0.999$  ( $\delta w^2, \delta b^2$ )

$\epsilon = 10^{-8}$

ADAM : "Adaptive moment Estimation."

### L9. LEARNING RATE DECAY:



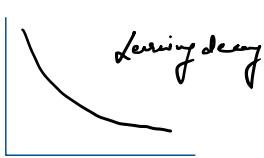
\* 1 epoch = 1 pass over the data.

$$\alpha = \left( \frac{1}{1 + \text{decay rate} * \text{epochnum}} \right) \alpha_0 \quad \xrightarrow{\text{some initial learning rate.}}$$

Epoch	$\alpha$
1	0.1
2	0.067
3	0.05
4	0.04

$$\alpha_0 = 0.2$$

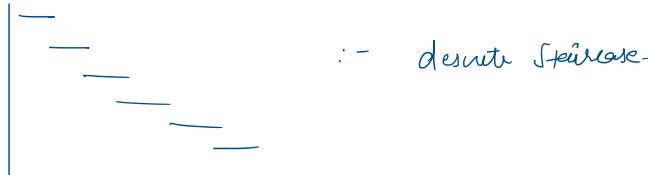
$$d\tau = 1$$



$$\alpha = \left[ (N < 1) \text{ epochnum.} \dots \alpha_0 \right] \vdash \text{exponential decay.}$$

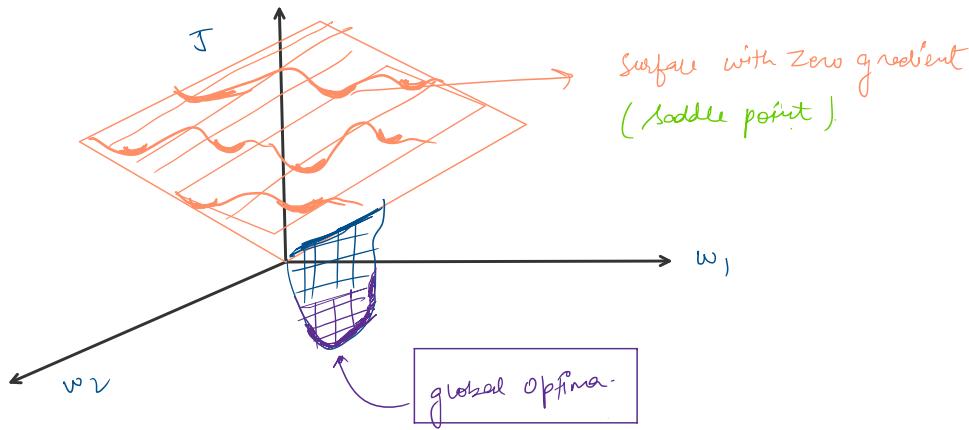
people also use:

$$\alpha = \frac{k}{\sqrt{\text{epoch}}} * \alpha_0$$



\* Manual decay ..

## L10: THE PROBLEM OF LOCAL OPTIMA



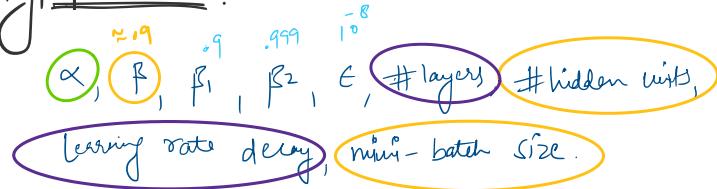
problem with plateaus.

- \* [ plateau is a region where the derivative is 0 for a very long time. ]
  - Unlikely to stuck in a bad local optima.
  - Plateaus can make learning slower.

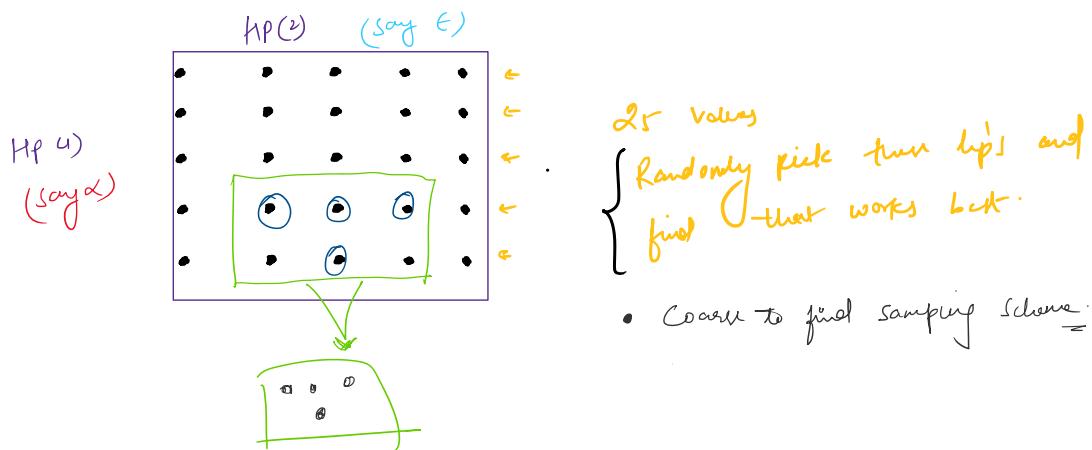
## Hyperparameter Tuning

### L1. Tuning process.

#### Hyperparameters



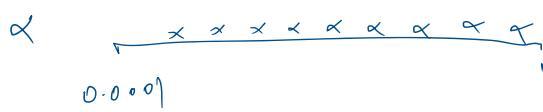
- \* - high priority.
- \* - medium priority.
- \* - low priority.
- \* - almost negligible (default values are used.)



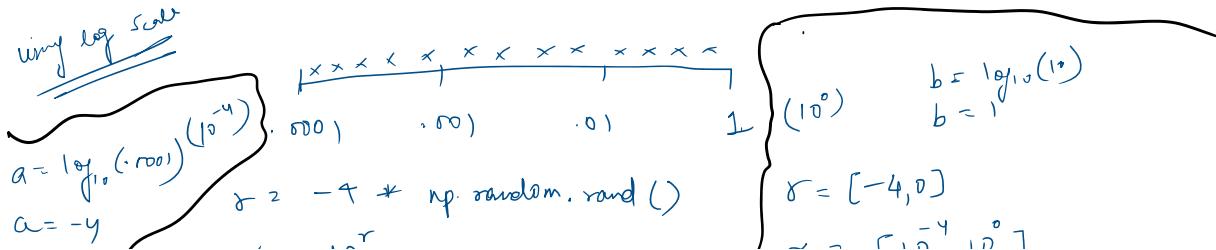
\* Use random Sampling and adequate search.

{optional} \* Coarse to fine search process.

### L2. Using An Appropriate Scale To Pick Hps.



90% Values  $\in [0.1, 1]$



$a = \log(1 - \text{rand})$   
 $a = -y$   
 $\alpha = 10^r$   
 Sample  $r$  uniformly on  $\in [a, b]$

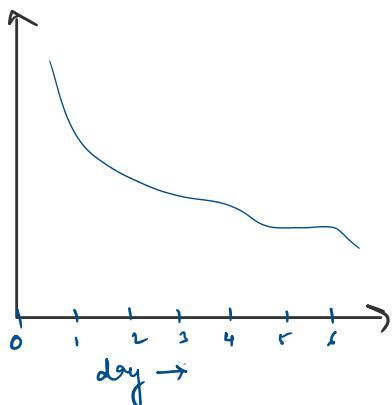
$\gamma = [-4, 0]$   
 $\alpha = [10^{-4}, 10^0]$

### # Hints for exponentially weighted avg.

$\beta = [0.9, \dots, 0.999]$   
 Linear Scale —  $\times$   
 $1 - \beta = [0.1, \dots, 0.001]$   
 1 → .1 → .01 → .001  
 ✓  
 very sensitive to small changes in  $\beta$ .

### L3: HPI tuning (Panda vs Caviar)

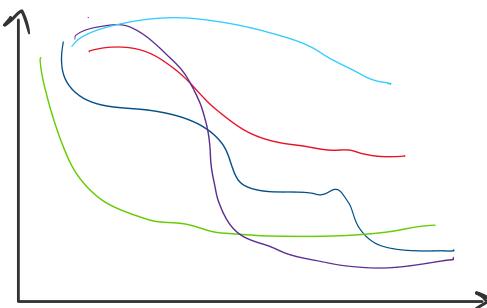
#### Babysitting One Model



Keep checking one day at a time.

#### Panda Approach

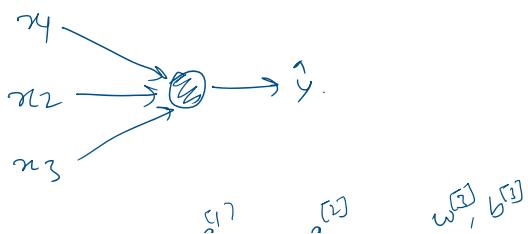
#### Training many models in Parallel



#### Caviar Approach

### Batch Normalization

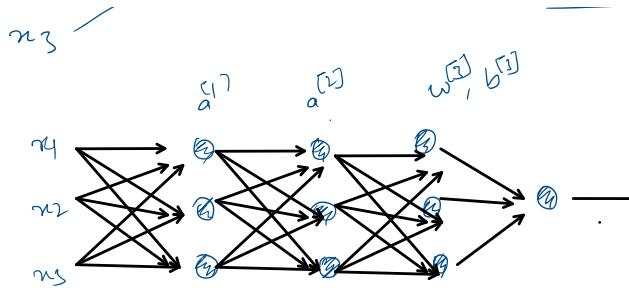
#### Normalizing Activations in a NW



$$x = x - \mu$$

$$\sigma^2 = \frac{1}{m} \sum_i x_i^2$$

$$x / \sigma^2$$



Can we normalize  $a^{[2]}$ ? to train  $w^{[3]}, b^{[3]}$  faster?

\* Normalize  $z^{[2]}$  or  $\tilde{z}^{[2]}$  instead of  $a^{[2]}$

$$\mu = \frac{1}{m} \sum z^{[i]}$$

$$\sigma^2 = \frac{1}{m} \sum (z^{[i]} - \mu)^2$$

$$z_{\text{normalize}}^{[i]} = \frac{z^{[i]} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$z$  has normal distribution and variance 1, but we want different dist.

$$\tilde{z}^{[i]} = \gamma z_{\text{norm}}^{[i]} + \beta$$

Learnable parameter.

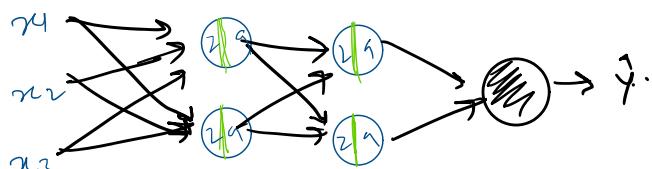
Standardize mean  
and Variance.

$$\gamma = \sqrt{\sigma^2 + \epsilon} \quad \text{and} \quad \beta = \mu$$

$$\tilde{z}^{[i]} = \gamma z^{[i]} + \beta$$

Use  $\tilde{z}^{[i]}$  instead of  $z^{[i]}$

## L2. Fitting Batch Norm to a N/W $\Rightarrow$



$$x \xrightarrow{w^{[0]}, b^{[0]}} z^{[0]} \xrightarrow{\gamma^{[0]}, \beta^{[0]}, \mu^{[0]}, \sigma^2} \tilde{z}^{[0]} \xrightarrow{a^{[0]} = g(\tilde{z}^{[0]})} \tilde{z}^{[1]} \dots$$

$$\text{Parameters} \Rightarrow w^{[0]}, b^{[0]}, \beta^{[0]}, \gamma^{[0]}, \mu^{[0]}, \sigma^2$$

$$\beta^{[i]} := \beta^{[i]} - \alpha \cdot d\beta^{[i]}$$

### Working with mini-batches.

$$x^{[1]} \xrightarrow{w^{[0]}, b^{[0]}} \tilde{z}^{[1]} \xrightarrow{\gamma^{[1]}, \beta^{[1]}, \mu^{[1]}, \sigma^2} g^{[1]}(\tilde{z}^{[1]}) \xrightarrow{a^{[1]}} \dots$$

$$\hat{x} = \overbrace{w^{[0]} b^{[0]}}^{\text{forward pass}} \xrightarrow{(f \cdot n)} \underbrace{0}_{\text{batch norm}} \xrightarrow{\gamma^{[0]} \beta^{[0]}} \hat{x}$$

$$\text{Parameters} = w^{[0]}, b^{[0]}, \beta^{[0]}, \gamma^{[0]}, \Rightarrow \hat{x} = w^{[0]} a + b^{[0]}$$

\* Batch Norm averages all  $\hat{x}^{[0]}$ , so adding/subtracting any const. doesn't affect the equation.

$$\hat{x}^{[0]} = w^{[0]} a^{[0]}$$

$$\hat{x}_{\text{norm.}}^{[0]} = \frac{\hat{x}^{[0]} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\tilde{x}^{[0]} = \gamma^{[0]} \hat{x}^{[0]} + \beta^{[0]}$$

$$\text{dim}_a := (n^{[0]}, 1)$$

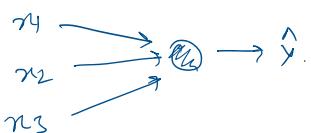
Implementing gradient descent.

for  $t=1$  to ... mini-Batch<sub>n</sub> :-

- ① Compute forward prop. on  $x^{[t]}$
- ② use batch norm to replace  $\hat{x}^{[0]} \Rightarrow \tilde{x}^{[0]}$
- ③ use backprop to calculate grads.
- ④ update parameters,  $\rightarrow w^{[0]} := w^{[0]} - \alpha \delta w^{[0]}$   
 $\rightarrow \beta^{[0]} := \beta^{[0]} - \alpha \delta \beta^{[0]}$   
 $\rightarrow \gamma^{[0]} := \gamma^{[0]} - \alpha \delta \gamma^{[0]}$

### L3. Why does batch norm works?

\* Leaving on shifting input dist.

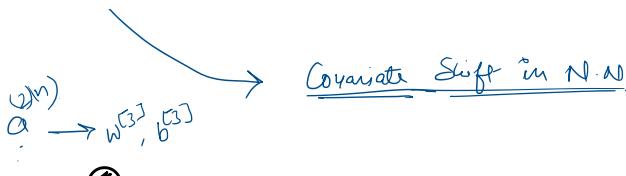


Trained on monochrome image

Trained on coloured image

"Covariate Shift"  $\rightarrow$

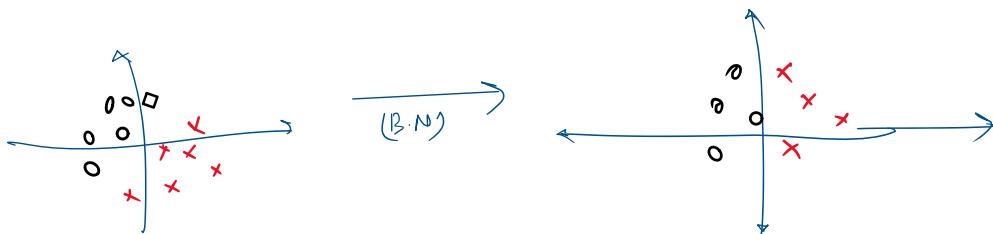
\* If you learn some  $x \rightarrow y$  mapping, if distribution of  $x$  changes then you need to retrain your learning algorithm.



~~Problem~~

parameters  $w, b$  before this layer change,  $a^{(2)}$  will also change.  
Suffering from a problem of "covariate shift".

- # Batch norm : No matter the distribution changes, mean and variance remains same.



## # Batch Norm as Regularization.

- ① Each mini-batch is scaled by mean/variance computed on just one batch.
- ② This adds some noise  $\tilde{z}^{(i)}$ , scaling process of  $z^{(i)} \rightarrow \tilde{z}^{(i)}$  becomes more noisier, similar to dropout.
- ③ This has slight regularization effect.

mini-batch size  $\propto \frac{1}{\text{noise}} \propto \text{regularization}$ .



"don't use batch norm as regularization  
only use for normalization."



## L4: Batch Norm at -last time.

We need different  $\mu$  and  $\sigma^2$ .  
|

$$\mu = \frac{1}{m} \sum z^{(i)}$$

$m$  :- ex: in mini-batches

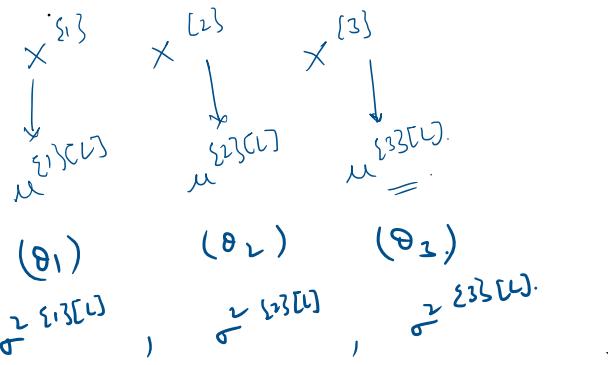
$$\sigma^2 = \frac{1}{m} \sum_i (z^{(i)} - \mu)^2$$

$$\rightarrow z^{(i)} - \mu$$

$$Z_{\text{norm}}^{(i)} = \frac{Z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\tilde{Z}^{(i)} = \gamma Z_{\text{norm}} + \beta$$

estimate using exponentially weighted avg. (across mini-batch)



new  
Znorm

$$\tilde{Z}_{\text{norm}} = \frac{\tilde{Z} - \mu}{\sqrt{\sigma^2 + \epsilon}}, \quad \tilde{Z} = \left[ \gamma Z_{\text{norm}} + \beta \right].$$

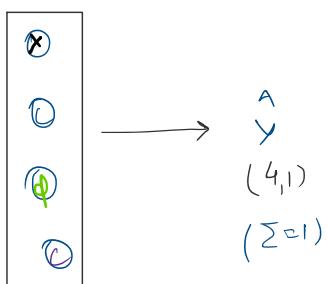
### Multi Class Classification

#### Softmax Regression.

Now recognize cats, dogs and chicks.

Class 1    Class 2    Class 3    -1

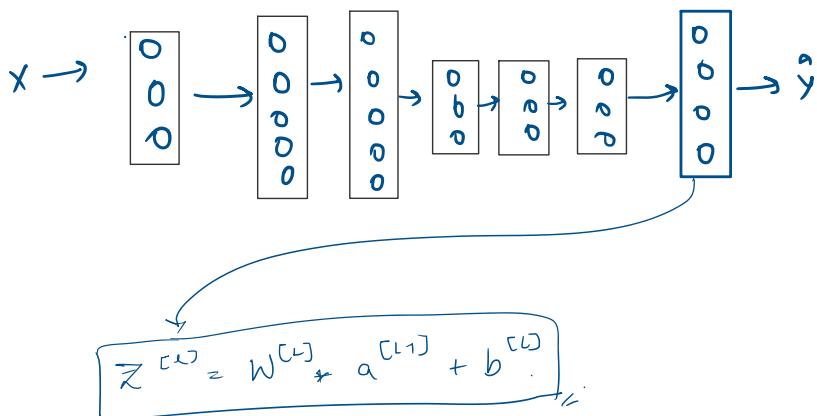
# classes : 4



- final layer -

Softmax Layer

## Softmax Layer



## Softmax Activation fn.

$$t = e^{(\tilde{z}^{[L]})}$$

$$(4,1) \quad a^{[L]} = \frac{e^{(\tilde{z}^{[L]})}}{\sum t^i}, \quad a_i^{[L]} = \frac{t^i}{\sum_{i=1}^4 t^i}$$

$$\tilde{z}^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix}$$

$$\sum_{i=1}^4 t^i = 176.3$$

$$a^{[L]} = \frac{t}{176.3} = \frac{1}{176.3} * \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix}$$

$$a^{[L]} = g^{(L)}(\tilde{z}^{[L]}).$$

$\uparrow_{4 \times 1}$

L2:

Training a softmax classifier:

$$\tilde{z}^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}, \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix}$$

$$a^{(v)} = g^{(v)}(z^{(v)}) = \begin{pmatrix} e^5 / \sum e \\ e^2 / \sum e \\ e^1 / \sum e \\ e^3 / \sum e \end{pmatrix} = \begin{pmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{pmatrix}$$

"hard max"

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow \text{Class with highest probability.}$$

- \* Softmax regression generalizes Logistic Regression to  $c$  classes.  
if  $c=2$ , Softmax Regression = Logistic Regression.

### Loss fx.

$$y^{(x_i)} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \xrightarrow{\text{Cat.}} \hat{y}^{(x_i)} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

$$L(\hat{y}, y) = - \sum_{j=1}^c y_j (\log \hat{y}_j)$$

$$(y_1 = y_3 = y_4 = 0) : 3 \text{ values in } L(\hat{y}, y) = 0$$

$$L(\hat{y}, y) = -y_2 \log \hat{y}_2$$

(minimize)  $\hat{y}_2 = 1 \Rightarrow -\log \hat{y}_2$

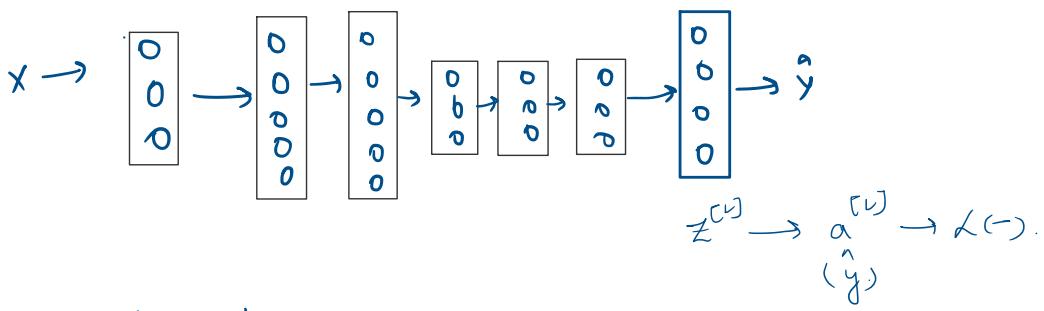
make  $\hat{y}$  as big

$$J(\omega^{(v)}, b^{(v)}) = 1/m \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

$$y = [y^{(1)} \ y^{(2)} \ y^{(3)} \dots \ y^{(m)}]_{(4 \times m)}$$

$$\hat{y} = [\hat{y}^{(1)} \ \hat{y}^{(2)} \ \hat{y}^{(3)} \dots \ \hat{y}^{(m)}]_{(4 \times m)}$$

### Gradient Descent with Softmax.



forward prop : ✓

backward prop :-

$$\delta z^{[L]} = \hat{y} - y \quad (\in \mathbb{R}^n) \quad (\leftarrow)$$

$\frac{\partial J}{\partial z^{[L]}}$

### Deep Learning Frameworks

- Caffe | Caffe2
- CNTK
- DL4J
- Keras
- Lasagne
- mx net
- PaddlePaddle
- PyTorch
- TensorFlow
- Theano